

Plots, Prediction and Testing for the Multivariate Linear Model

David J. Olive*

Southern Illinois University

December 18, 2013

Abstract

To use the multivariate linear model, reliable prediction regions and tests of hypotheses are needed. This paper provides useful prediction regions, gives F approximations to the widely used Wilks' Λ , Pillai's trace, and Hotelling Lawley trace test statistics, and gives plots to check goodness and lack of fit, to check for outliers and influential cases, and to check whether the error distribution is multivariate normal or from some other elliptically contoured distribution. Some of the plots and prediction regions can be extended to more general multivariate regression models.

KEY WORDS: MANOVA; multivariate linear regression; multivariate regression; prediction regions; seemingly unrelated regressions.

*David J. Olive is Associate Professor, Department of Mathematics, Southern Illinois University, Mailcode 4408, Carbondale, IL 62901-4408, USA. E-mail address: dolive@siu.edu.

1 INTRODUCTION

For the multivariate linear model, this paper provides useful prediction regions and plots, and gives F approximations to the widely used Wilks' Λ , Pillai's trace, and Hotelling Lawley trace test statistics. First notation is given for the multivariate linear model.

The *multivariate linear model* $\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \boldsymbol{\epsilon}_i$ for $i = 1, \dots, n$ has $m \geq 2$ response variables Y_1, \dots, Y_m and p predictor variables x_1, x_2, \dots, x_p . Multivariate linear regression and MANOVA models are special cases. The i th case is $(\mathbf{x}_i^T, \mathbf{y}_i^T) = (x_{i1}, x_{i2}, \dots, x_{ip}, Y_{i1}, \dots, Y_{im})$. If a constant $x_{i1} = 1$ is in the model, then x_{i1} could be omitted from the case. The model is written in matrix form as $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$ where the matrices are defined below. The model has $E(\boldsymbol{\epsilon}_k) = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\epsilon}_k) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} = (\sigma_{ij})$ for $k = 1, \dots, n$. Also $E(\mathbf{e}_i) = \mathbf{0}$ while $\text{Cov}(\mathbf{e}_i, \mathbf{e}_j) = \sigma_{ij} \mathbf{I}_n$ for $i, j = 1, \dots, m$ where \mathbf{I}_n is the $n \times n$ identity matrix and \mathbf{e}_i is defined below. Then the $p \times m$ coefficient matrix $\mathbf{B} = \begin{bmatrix} \boldsymbol{\beta}_1 & \boldsymbol{\beta}_2 & \dots & \boldsymbol{\beta}_m \end{bmatrix}$ and the $m \times m$ covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$ are to be estimated, and $E(\mathbf{Z}) = \mathbf{X}\mathbf{B}$ while $E(Y_{ij}) = \mathbf{x}_i^T \boldsymbol{\beta}_j$.

The $n \times m$ matrix of response variables and $n \times m$ matrix of errors are

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Y}_1 & \mathbf{Y}_2 & \dots & \mathbf{Y}_m \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_n^T \end{bmatrix} \quad \text{and} \quad \mathbf{E} = \begin{bmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \dots & \mathbf{e}_m \end{bmatrix} = \begin{bmatrix} \boldsymbol{\epsilon}_1^T \\ \vdots \\ \boldsymbol{\epsilon}_n^T \end{bmatrix},$$

while the $n \times p$ design matrix of predictor variables is \mathbf{X} .

Least squares is the classical method for fitting the multivariate linear model. The *least squares estimators* are $\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Z} = \begin{bmatrix} \hat{\boldsymbol{\beta}}_1 & \hat{\boldsymbol{\beta}}_2 & \dots & \hat{\boldsymbol{\beta}}_m \end{bmatrix}$. The matrix of *predicted values* or *fitted values* $\hat{\mathbf{Z}} = \mathbf{X}\hat{\mathbf{B}} = \begin{bmatrix} \hat{\mathbf{Y}}_1 & \hat{\mathbf{Y}}_2 & \dots & \hat{\mathbf{Y}}_m \end{bmatrix}$. The matrix of *residuals* $\hat{\mathbf{E}} = \mathbf{Z} - \hat{\mathbf{Z}} = \mathbf{Z} - \mathbf{X}\hat{\mathbf{B}} = \begin{bmatrix} \hat{\mathbf{r}}_1 & \hat{\mathbf{r}}_2 & \dots & \hat{\mathbf{r}}_m \end{bmatrix}$. These quantities can be found

from the m multiple linear regressions of Y_j on the predictors: $\hat{\boldsymbol{\beta}}_j = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}_j$, $\hat{\mathbf{Y}}_j = \mathbf{X} \hat{\boldsymbol{\beta}}_j$ and $\hat{\mathbf{r}}_j = \mathbf{Y}_j - \hat{\mathbf{Y}}_j$ for $j = 1, \dots, m$. Hence $\hat{\epsilon}_{i,j} = Y_{i,j} - \hat{Y}_{i,j}$ where $\hat{\mathbf{Y}}_j = (\hat{Y}_{1,j}, \dots, \hat{Y}_{n,j})^T$. Finally,

$$\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},d} = \frac{(\mathbf{Z} - \hat{\mathbf{Z}})^T (\mathbf{Z} - \hat{\mathbf{Z}})}{n-d} = \frac{(\mathbf{Z} - \mathbf{X} \hat{\mathbf{B}})^T (\mathbf{Z} - \mathbf{X} \hat{\mathbf{B}})}{n-d} = \frac{\hat{\mathbf{E}}^T \hat{\mathbf{E}}}{n-d} = \frac{1}{n-d} \sum_{i=1}^n \hat{\boldsymbol{\epsilon}}_i \hat{\boldsymbol{\epsilon}}_i^T.$$

The choices $d = 0$ and $d = p$ are common. If $d = 1$, then $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},d=1} = \mathbf{S}_r$, the sample covariance matrix of the residual vectors $\hat{\boldsymbol{\epsilon}}_i$ since the sample mean of the $\hat{\boldsymbol{\epsilon}}_i$ is $\mathbf{0}$. Let $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},p}$ be the unbiased estimator of $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$.

The $\boldsymbol{\epsilon}_i$ are assumed to be iid. Some important joint distributions for $\boldsymbol{\epsilon}$ are completely specified by an $m \times 1$ population *location* vector $\boldsymbol{\mu}$ and an $m \times m$ symmetric positive definite population *dispersion* matrix $\boldsymbol{\Sigma}$. An important model is the elliptically contoured $EC_m(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$ distribution with probability density function

$$f(\mathbf{z}) = k_m |\boldsymbol{\Sigma}|^{-1/2} g[(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu})]$$

where $k_m > 0$ is some constant and g is some known function. The multivariate normal (MVN) $N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ distribution is a special case.

Some additional notation will be useful. Assume that $\mathbf{x}_1, \dots, \mathbf{x}_n$ are iid from a multivariate distribution. The classical estimator $(\bar{\mathbf{x}}, \mathbf{S})$ of multivariate location and dispersion is the sample mean and sample covariance matrix where

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \text{and} \quad \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T. \quad (1)$$

Let the $p \times 1$ column vector T be a multivariate location estimator, and let the $p \times p$ symmetric positive definite matrix \mathbf{C} be a dispersion estimator. Then the i th *squared*

sample Mahalanobis distance is the scalar

$$D_i^2 = D_i^2(T, \mathbf{C}) = (\mathbf{x}_i - T)^T \mathbf{C}^{-1} (\mathbf{x}_i - T) \quad (2)$$

for each observation \mathbf{x}_i . Notice that the Euclidean distance of \mathbf{x}_i from the estimate of center T is $D_i(T, \mathbf{I}_p)$. The classical Mahalanobis distance uses $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$. Following Johnson (1987, pp. 107-108), the population squared Mahalanobis distance

$$U \equiv D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}), \quad (3)$$

and for elliptically contoured distributions, U has probability density function (pdf)

$$h(u) = \frac{\pi^{p/2}}{\Gamma(p/2)} k_p u^{p/2-1} g(u). \quad (4)$$

The classical large sample $100(1-\delta)\%$ prediction region for a future value \mathbf{x}_f given iid data $\mathbf{x}_1, \dots, \mathbf{x}_n$ is $\{\mathbf{x} : D_{\mathbf{x}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq \chi_{p,1-\delta}^2\}$, while for multivariate linear regression, the classical large sample $100(1-\delta)\%$ prediction region for a future value \mathbf{y}_f given \mathbf{x}_f and past data $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$ is $\{\mathbf{y} : D_{\mathbf{y}}^2(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_{\epsilon}) \leq \chi_{m,1-\delta}^2\}$. See Johnson and Wichern (1988, pp. 134, 151, 312). By (4), these regions may work for multivariate normal data, but otherwise tend to have undercoverage. Olive (2013) replaced $\chi_{p,1-\delta}^2$ by the order statistic $D_{(U_n)}^2$ where U_n decreases to $\lceil n(1-\delta) \rceil$. Section 3 will use a similar technique to develop possibly the first practical large sample prediction region for the multivariate linear model with unknown error distribution.

Following regression graphics notation, suppress the subscript indicating case for the response variable and vector of predictor variables. Suppose $m = 1$ and the response variable Y is conditionally independent of the vector of predictors \mathbf{x} given the *sufficient predictor* $SP = h(\mathbf{x})$, where the d -dimensional function $h : \mathcal{R}^p \rightarrow \mathcal{R}^d$. The *estimated*

sufficient predictor $\text{ESP} = \hat{h}(\mathbf{x})$. If $d = 1$, then a *response plot* is a plot of ESP versus Y . Cook and Weisberg (1999a, p. 411) define a *sufficient summary plot* to be a plot that contains all of the sample regression information about the conditional distribution of the response given the predictors. Hence a plot of $h(\mathbf{x})$ versus Y is a sufficient summary plot while a response plot of $\hat{h}(\mathbf{x})$ versus Y is an estimated sufficient summary plot. As an analogy, a sufficient statistic contains all of the sample information for estimating a parameter θ , and want the dimension of the sufficient statistic as small as possible. Similarly want d as small as possible for the greatest amount of dimension reduction.

The Rousseeuw and Van Driessen (1999) DD plot is a plot of classical Mahalanobis distances versus robust Mahalanobis distances. Simulations and results from Olive (2002) suggest that as $n \rightarrow \infty$, the plotted points in the DD plot of the least squares residual vectors $\hat{\boldsymbol{\epsilon}}_i$ will cluster tightly about the identity line with unit slope and zero intercept if the $\boldsymbol{\epsilon}_i$ are iid from a multivariate normal $N_m(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}})$ distribution, and cluster tightly about some line through the origin with slope greater than one for a large class of elliptically contoured distributions.

Section 2 suggests making the m response plots, the m residual plots and the DD plot of the residual vectors. For the multivariate linear model, these plots form the minimal collection of 2-dimensional plots for visualizing the conditional distribution of $\mathbf{y}|\mathbf{x}$. For this model, $\hat{h}_j(\mathbf{x}) = \mathbf{x}^T \hat{\boldsymbol{\beta}}_j = \hat{Y}_j$ for $j = 1, \dots, m = d$. Since Y_j depends on \mathbf{x} through $\mathbf{x}^T \boldsymbol{\beta}_j$, the conditional distribution of $Y_j|\mathbf{x}$ can be visualized with the response plot of $\mathbf{x}^T \hat{\boldsymbol{\beta}}_j$ versus Y_j . The response and residual plots will be used to check the multivariate linear model for linearity, influential cases, and outliers. The DD plot of the residuals is used to check the error distribution, to check for outliers, and to display the prediction

region developed in section 3.

Kakizawa (2009) examines testing for the multivariate linear regression model, showing that the Wilks, Pillai, and Hotelling Lawley test statistics perform well asymptotically for a large class of zero mean error distributions. Section 4 reviews these results and shows that the Hotelling Lawley test statistic is closely related to the partial F statistic for multiple linear regression. Section 5 gives some examples and simulations.

2 Plots for the Multivariate Linear Model

This section suggests using residual plots, response plots, and the DD plot to examine the multivariate linear model. The residual plots are often used to check for lack of fit of the multivariate linear model. The response plots are used to check linearity and to detect influential cases and outliers. The response and residual plots are used exactly as in the $m = 1$ case corresponding to multiple linear regression and experimental design models. See Olive and Hawkins (2005) and Cook and Weisberg (1999a, p. 432; 1999b)

The DD plot of the residual vectors $\hat{\boldsymbol{\epsilon}}_i$ is used to check the error distribution, to detect outliers and to display the nonparametric prediction region developed in section 3. The DD plot suggests that the error distribution is elliptically contoured if the plotted points cluster tightly about a line through the origin as $n \rightarrow \infty$. The plot suggests that the error distribution is multivariate normal if the line is the identity line. If n is large and the plotted points do not cluster tightly about a line through the origin, then the error distribution may not be elliptically contoured. Make a DD plot of the continuous predictor variables to check for \boldsymbol{x} -outliers. These applications of the DD plot for iid

multivariate data are discussed in Olive (2002, 2013). See Examples 1 and 2.

Make the m response and residual plots for the multivariate linear model. A *response plot* for the j th response variable is a plot of the fitted values \hat{Y}_{ij} versus the response Y_{ij} where $i = 1, \dots, n$. The identity line is added to the plot as a visual aid. A *residual plot* corresponding to the j th response variable is a plot of \hat{Y}_{ij} versus r_{ij} . In a response plot, the vertical deviations from the identity line are the residuals $r_{ij} = Y_{ij} - \hat{Y}_{ij}$. Assume that the \hat{Y}_{ij} take on many values (this assumption is not met by the one way MANOVA model, but is met by the multivariate linear regression model and by some MANOVA models). Suppose the model is good, the error distribution is not highly skewed, and $n \geq 10p$. Then the plotted points should cluster about the identity line in each of the m response plots. If outliers are present or if the plot is not linear, then the current model or data need to be changed or corrected. See Example 1 in section 5.

The response and residual plots for the one way MANOVA model need some notation, and it is useful to use three subscripts. Suppose there are independent random samples from p different populations (treatments), and n_i cases are randomly assigned to p treatment groups with $n = \sum_{i=1}^p n_i$. Assume that m response variables $\mathbf{y}_{ij} = (Y_{ij1}, \dots, Y_{ijm})^T$ are measured for the i th treatment. Hence $i = 1, \dots, p$ and $j = 1, \dots, n_i$. The Y_{ijk} follow different one way ANOVA models for $k = 1, \dots, m$. Assume $E(\mathbf{y}_{ij}) = \boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{im})^T$ and $\text{Cov}(\mathbf{y}_{ij}) = \boldsymbol{\Sigma}_\epsilon$. Hence the p treatments have possibly different mean vectors $\boldsymbol{\mu}_i$, but common covariance matrix $\boldsymbol{\Sigma}_\epsilon$.

Then for the k th response variable, the *response plot* is a plot of $\hat{Y}_{ijk} \equiv \hat{\mu}_{iok}$ versus Y_{ijk} and the *residual plot* is a plot of $\hat{Y}_{ijk} \equiv \hat{\mu}_{iok}$ versus r_{ijk} where the $\hat{\mu}_{iok}$ are the sample means of the p treatments for the k th response variable. Add the identity line to the

response plot and $r = 0$ line to the residual plot as visual aids. The points in the response plot scatter about the identity line and the points in the residual plot scatter about the $r = 0$ line. The response plot consists of p dot plots, one for each value of $\hat{\mu}_{iok}$. The dot plot corresponding to $\hat{\mu}_{iok}$ is the dot plot of $Y_{i,1,k}, \dots, Y_{i,n_i,k}$. Similarly, the residual plot consists of p dot plots, and the plot corresponding to $\hat{\mu}_{iok}$ is the dot plot of $r_{i,1,k}, \dots, r_{i,n_i,k}$. Assuming the $n_i \geq 10$, the p dot plots should have roughly the same amount of spread in both the response and residual plots. Again m response and residual plots are made, one for each response variable. See Example 2 in section 5.

3 Asymptotically Optimal Prediction Regions

In this section, we will consider a more general multivariate regression model, and then consider the multivariate linear model as a special case. Given n cases of training or past data $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$ and a vector of predictors \mathbf{x}_f , suppose it is desired to predict a future vector \mathbf{y}_f . Then a large sample $(1 - \delta)100\%$ prediction region is a set \mathcal{A}_n such that $P(\mathbf{y}_f \in \mathcal{A}_n) \rightarrow 1 - \delta$ as $n \rightarrow \infty$, and is asymptotically optimal if the volume of the region converges in probability to the volume of the population minimum volume covering region. The following technical theorem will be needed to prove Theorem 2, and the proof is in the appendix.

Theorem 1. Let $a_j > 0$ and assume that $(\hat{\boldsymbol{\mu}}_{j,n}, \hat{\boldsymbol{\Sigma}}_{j,n})$ are consistent estimators of $(\boldsymbol{\mu}, a_j \boldsymbol{\Sigma})$ for $j = 1, 2$.

a) $D_{\mathbf{x}}^2(\hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) - \frac{1}{a_j} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = o_P(1)$.

b) Let $0 < \delta \leq 0.5$. If $(\hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) - (\boldsymbol{\mu}, a_j \boldsymbol{\Sigma}) = O_p(n^{-\delta})$ and $a_j \hat{\boldsymbol{\Sigma}}_j^{-1} - \boldsymbol{\Sigma}^{-1} = O_P(n^{-\delta})$,

then

$$D_{\mathbf{x}}^2(\hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) - \frac{1}{a_j} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = O_P(n^{-\delta}).$$

Now suppose a prediction region for an $m \times 1$ random vector \mathbf{y}_f given a vector of predictors \mathbf{x}_f is desired for the multivariate linear model. If we had many cases $\mathbf{z}_i = \mathbf{B}^T \mathbf{x}_f + \boldsymbol{\epsilon}_i$, then we could make a prediction region for \mathbf{z}_i using a multivariate location and dispersion model prediction region for m variables. Instead, Theorem 2 will use a prediction region on the pseudodata $\hat{\mathbf{z}}_i = \hat{\mathbf{B}}^T \mathbf{x}_f + \hat{\boldsymbol{\epsilon}}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$ for $i = 1, \dots, n$. This takes the data cloud of the n residual vectors $\hat{\boldsymbol{\epsilon}}_i$ and centers the cloud at $\hat{\mathbf{y}}_f$. Note that $\hat{\mathbf{z}}_i = (\mathbf{B} - \mathbf{B} + \hat{\mathbf{B}})^T \mathbf{x}_f + (\boldsymbol{\epsilon}_i - \boldsymbol{\epsilon}_i + \hat{\boldsymbol{\epsilon}}_i) = \mathbf{z}_i + (\hat{\mathbf{B}} - \mathbf{B})^T \mathbf{x}_f + \hat{\boldsymbol{\epsilon}}_i - \boldsymbol{\epsilon}_i = \mathbf{z}_i + (\hat{\mathbf{B}} - \mathbf{B})^T \mathbf{x}_f - (\hat{\mathbf{B}} - \mathbf{B})^T \mathbf{x}_i = \mathbf{z}_i + O_P(n^{-1/2})$.

If the $\boldsymbol{\epsilon}_i$ are iid from an $EC_m(\mathbf{0}, \boldsymbol{\Sigma}, g)$ distribution with continuous decreasing g and nonsingular covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} = c\boldsymbol{\Sigma}$ for some constant $c > 0$, then the population asymptotically optimal prediction region is $\{\mathbf{y} : D_{\mathbf{y}}(\mathbf{B}^T \mathbf{x}_f, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}\}$ where $P(D_{\mathbf{y}}(\mathbf{B}^T \mathbf{x}_f, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}) = 1 - \delta$. For example, if the iid $\boldsymbol{\epsilon}_i \sim N_m(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}})$, then $D_{1-\delta} = \sqrt{\chi_{m,1-\delta}^2}$. If the error distribution is not elliptically contoured, then the above region still has $100(1 - \delta)\%$ coverage, but prediction regions with smaller volume may exist.

A natural way to make a large sample prediction region is to estimate the target population minimum volume covering region, but for moderate samples and many error distributions, the natural estimator that covers $\lceil n(1 - \delta) \rceil$ of the cases tends to have undercoverage as high as $\min(0.05, \delta/2)$. This empirical result is not too surprising since it is well known that the performance of a prediction region on the training data is

superior to the performance on future data, due in part to the unknown variability of the estimator. To compensate for the undercoverage, let q_n be as in Theorem 2.

Theorem 2. Suppose $\mathbf{y}_i = E(\mathbf{y}_i) + \boldsymbol{\epsilon}_i = \hat{\mathbf{y}}_i + \hat{\boldsymbol{\epsilon}}_i$ where $\text{Cov}(\boldsymbol{\epsilon}_i) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} > 0$, and where the zero mean $\boldsymbol{\epsilon}_f$ and the $\boldsymbol{\epsilon}_i$ are iid for $i = 1, \dots, n$. Given \mathbf{x}_f , suppose the fitted model produces $\hat{\mathbf{y}}_f$ and nonsingular $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}$. Let $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$ and

$$D_i^2(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}) = (\hat{\mathbf{z}}_i - \hat{\mathbf{y}}_f)^T \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} (\hat{\mathbf{z}}_i - \hat{\mathbf{y}}_f)$$

for $i = 1, \dots, n$. Let $q_n = \min(1 - \delta + 0.05, 1 - \delta + m/n)$ for $\delta > 0.1$ and

$$q_n = \min(1 - \delta/2, 1 - \delta + 10\delta m/n), \quad \text{otherwise.}$$

If $q_n < 1 - \delta + 0.001$, set $q_n = 1 - \delta$. Let $0 < \delta < 1$ and $h = D_{(U_n)}$ where $D_{(U_n)}$ is the q_n th sample quantile of the D_i . Let the nominal $100(1 - \delta)\%$ prediction region for \mathbf{y}_f be given by $\{\mathbf{z} : (\mathbf{z} - \hat{\mathbf{y}}_f)^T \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} (\mathbf{z} - \hat{\mathbf{y}}_f) \leq D_{(U_n)}^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}) \leq D_{(U_n)}^2\} =$

$$\{\mathbf{z} : D_{\mathbf{z}}(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}) \leq D_{(U_n)}\}. \quad (5)$$

a) Consider the n prediction regions for the data where $(\mathbf{y}_{f,i}, \mathbf{x}_{f,i}) = (\mathbf{y}_i, \mathbf{x}_i)$ for $i = 1, \dots, n$. If the order statistic $D_{(U_n)}$ is unique, then U_n of the n prediction regions contain \mathbf{y}_i where $U_n/n \rightarrow 1 - \delta$ as $n \rightarrow \infty$.

b) If $(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}})$ is a consistent estimator of $(E(\mathbf{y}_f), \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}})$, then (5) is a large sample $100(1 - \delta)\%$ prediction region for \mathbf{y}_f .

c) If $(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}})$ is a consistent estimator of $(E(\mathbf{y}_f), \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}})$, and the $\boldsymbol{\epsilon}_i$ come from an elliptically contoured distribution such that the highest density region is $\{\mathbf{z} : D_{\mathbf{z}}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}\}$, then the prediction region (5) is asymptotically optimal.

The proof of Theorem 2 is in the appendix. Notice that for the data $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$, if $\hat{\Sigma}_{\epsilon}^{-1}$ exists, then $100q_n\%$ of the n cases are in their corresponding prediction region, and $q_n \rightarrow 1 - \delta$ even if $(\hat{\mathbf{y}}_i, \hat{\Sigma}_{\epsilon})$ is not a good estimator. Hence the coverage q_n of the training data is robust to model assumptions. Of course the volume of the prediction region could be large if a poor estimator $(\hat{\mathbf{y}}_i, \hat{\Sigma}_{\epsilon})$ is used or if the ϵ_i do not come from an elliptically contoured distribution. The response, residual and DD plots can be used to check model assumptions.

Also want $n \geq 10 \max(p, m)$ for least squares estimators. If n is too small, then multivariate data is sparse and the covering ellipsoid for the training data may be far too small for future data, resulting in severe undercoverage. Coverage can also be arbitrarily bad if there is extrapolation or if $(\mathbf{x}_f, \mathbf{y}_f)$ comes from a different population than that of the data. Also notice that $q_n = 1 - \delta/2$ or $q_n = 1 - \delta + 0.05$ for $n \leq 20p$. At the training data, the coverage $q_n \geq 1 - \delta$, and q_n converges to the nominal coverage $1 - \delta$ as $n \rightarrow \infty$. Suppose $n \leq 20p$. Then the nominal 95% prediction region uses $q_n = 0.975$ while the nominal 50% prediction region uses $q_n = 0.55$. Prediction distributions depend both on the error distribution and on the variability of the estimator $(\hat{\mathbf{y}}_f, \hat{\Sigma}_{\epsilon})$. This variability is typically unknown but converges to 0 as $n \rightarrow \infty$. Also, residuals tend to underestimate errors for small n . For moderate n , ignoring estimator variability and using $q_n = 1 - \delta$ resulted in undercoverage as high as $\min(0.05, \delta/2)$. Letting the “coverage” q_n decrease to the nominal coverage $1 - \delta$ inflates the volume of the prediction region for small n , compensating for the unknown variability of $(\hat{\mathbf{y}}_f, \hat{\Sigma}_{\epsilon})$.

Prediction region (5) can be used for multivariate linear regression using least squares or the Su and Cook (2012) inner envelopes estimator. The region can also be used for

the seemingly unrelated regressions model. Theorem 3 relates prediction region (5) to a prediction region originally created for iid multivariate data.

Theorem 3. For multivariate linear regression, when least squares is used to compute $\hat{\mathbf{y}}_f$, \mathbf{S}_r and the pseudodata $\hat{\mathbf{z}}_i$, prediction region (5) is the Olive (2013) nonparametric prediction region applied to the $\hat{\mathbf{z}}_i$.

The proof of Theorem 3 is in the appendix. Olive (2013) derived prediction regions for a future observation \mathbf{x}_f given n iid $p \times 1$ random vectors \mathbf{x}_i . Suppose (T, \mathbf{C}) is an estimator of multivariate location and dispersion $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ such as the classical estimator $(\bar{\mathbf{x}}, \mathbf{S})$ given by (1). For $h > 0$, consider the hyperellipsoid

$$\{\mathbf{z} : (\mathbf{z} - T)^T \mathbf{C}^{-1} (\mathbf{z} - T) \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2 \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}} \leq h\}. \quad (6)$$

A future observation \mathbf{x}_f is in region (6) if $D_{\mathbf{x}_f} \leq h$. If (T, \mathbf{C}) is a consistent estimator of $(\boldsymbol{\mu}, d\boldsymbol{\Sigma})$, then (6) is a large sample $(1 - \delta)100\%$ prediction region if $h = D_{(U_n)}$ where $D_{(U_n)}$ is the q_n th sample quantile of the D_i with p replacing m . The classical parametric multivariate normal large sample prediction region uses $D_{\mathbf{x}_f} \equiv MD_{\mathbf{x}_f} \leq \sqrt{\chi_{p,1-\delta}^2}$.

The Olive and Hawkins (2010) RMVN estimator $(T_{RMVN}, \mathbf{C}_{RMVN})$ is an easily computed \sqrt{n} consistent estimator of $(\boldsymbol{\mu}, c\boldsymbol{\Sigma})$ under regularity conditions (E1) that include a large class of elliptically contoured distributions, and $c = 1$ for the $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ distribution. The RMVN estimator also gives useful estimates of $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ for $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ data even when certain types of outliers are present, and will be the robust estimator used in the DD plots. Also see Zhang, Olive and Ye (2012).

The nonparametric region uses the classical estimator $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ and $h = D_{(U_n)}$. The semiparametric region uses $(T, \mathbf{C}) = (T_{RMVN}, \mathbf{C}_{RMVN})$ and $h = D_{(U_n)}$.

The parametric MVN region uses $(T, \mathbf{C}) = (T_{RMVN}, \mathbf{C}_{RMVN})$ and $h^2 = \chi_{p, q_n}^2$ where $P(W \leq \chi_{p, q_n}^2) = q_n$ if $W \sim \chi_p^2$. All three regions are asymptotically optimal for $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ distributions with nonsingular $\boldsymbol{\Sigma}$. Olive (2013) shows that the first two regions are asymptotically optimal under (E1). For distributions with nonsingular covariance matrix $c_X \boldsymbol{\Sigma}$, the nonparametric region is a large sample $(1-\delta)100\%$ prediction region, but regions with smaller volume may exist.

The nonparametric prediction region for multivariate linear regression of Theorem 3 uses $(T, \mathbf{C}) = (\hat{\mathbf{y}}_f, \mathbf{S}_r)$ in (5), and has simple geometry. Let R_r be the nonparametric prediction region applied to the residuals $\hat{\epsilon}_i$. Then R_r is a hyperellipsoid with center $\mathbf{0}$, and the nonparametric prediction region is the hyperellipsoid R_r translated to have center $\hat{\mathbf{y}}_f$. See examples in section 5.

4 Testing Hypotheses

Consider testing a linear hypothesis $H_0 : \mathbf{L}\mathbf{B} = \mathbf{0}$ versus $H_1 : \mathbf{L}\mathbf{B} \neq \mathbf{0}$ where \mathbf{L} is a full rank $r \times p$ matrix. Let $\mathbf{H} = \hat{\mathbf{B}}^T \mathbf{L}^T [\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} \mathbf{L} \hat{\mathbf{B}}$. Let the error or residual sum of squares and cross products matrix be

$$\mathbf{W}_e = \hat{\mathbf{E}}^T \hat{\mathbf{E}} = (\mathbf{Z} - \hat{\mathbf{Z}})^T (\mathbf{Z} - \hat{\mathbf{Z}}) = \mathbf{Z}^T \mathbf{Z} - \mathbf{Z}^T \mathbf{X} \hat{\mathbf{B}} = \mathbf{Z}^T [\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{Z}.$$

Then $\mathbf{W}_e / (n - p) = \hat{\boldsymbol{\Sigma}} \boldsymbol{\epsilon}$. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ be the ordered eigenvalues of $\mathbf{W}_e^{-1} \mathbf{H}$.

Then there are four commonly used test statistics.

The Roy's maximum root statistic is $\lambda_{max}(\mathbf{L}) = \lambda_1$.

The Wilks' Λ statistic is $\Lambda(\mathbf{L}) = |(\mathbf{H} + \mathbf{W}_e)^{-1} \mathbf{W}_e| = |\mathbf{W}_e^{-1} \mathbf{H} + \mathbf{I}|^{-1} = \prod_{i=1}^m (1 + \lambda_i)^{-1}$.

The Pillai's trace statistic is $V(\mathbf{L}) = \text{tr}[(\mathbf{H} + \mathbf{W}_e)^{-1}\mathbf{H}] = \sum_{i=1}^m \frac{\lambda_i}{1 + \lambda_i}$.

The Hotelling-Lawley trace statistic is $U(\mathbf{L}) = \text{tr}[\mathbf{W}_e^{-1}\mathbf{H}] = \sum_{i=1}^m \lambda_i$.

Theorem 4. The Hotelling-Lawley trace statistic

$$U(\mathbf{L}) = \frac{1}{n-p} [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\Sigma}_\epsilon^{-1} \otimes (\mathbf{L}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]. \quad (7)$$

The proof of Theorem 4 is in the appendix. Some notation is useful to show (7) and to show that $(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$ under mild conditions if H_0 is true. Following Henderson and Searle (1979), let matrix $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_p]$. Then the vec operator stacks the columns of \mathbf{A} on top of one another, and $\mathbf{A} \otimes \mathbf{B}$ is the Kronecker product of \mathbf{A} and \mathbf{B} . An important fact is that if \mathbf{A} and \mathbf{B} are nonsingular square matrices, then $[\mathbf{A} \otimes \mathbf{B}]^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$. The following assumption is important.

Assumption D1: Let h_i be the i th diagonal element of $\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$. Assume $\max_{1 \leq i \leq n} h_i \rightarrow 0$ as $n \rightarrow \infty$, assume that the zero mean iid errors have finite fourth moments, and assume that $\frac{1}{n}\mathbf{X}^T\mathbf{X} \xrightarrow{P} \mathbf{W}^{-1}$.

Then for the least squares estimator, Su and Cook (2012) show that if assumption D1 holds, then $\hat{\Sigma}_\epsilon$ is \sqrt{n} consistent and $\sqrt{n} \text{vec}(\hat{\mathbf{B}} - \mathbf{B}) \xrightarrow{D} N_{pm}(\mathbf{0}, \Sigma_\epsilon \otimes \mathbf{W})$.

Theorem 5. If assumption D1 holds and if H_0 is true, then $(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$.

The proof of Theorem 5 is in the appendix. Kakizawa (2009) shows, under stronger assumptions than Theorem 5, that for a large class of iid error distributions, the following test statistics have the same χ_{rm}^2 limiting distribution when H_0 is true, and the same noncentral $\chi_{rm}^2(\omega^2)$ limiting distribution with noncentrality parameter ω^2 when H_0 is false under a local alternative. Hence the three tests are robust to the assumption of normality. The limiting null distribution is well known when the zero mean errors

are iid from a multivariate normal distribution. See Khattree and Naik (1999, p. 68): $(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$, $(n-p)V(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$, and $-[n-p-0.5(m-r+3)] \log(\Lambda(\mathbf{L})) \xrightarrow{D} \chi_{rm}^2$. Results from Kshirsagar (1972, p. 301) suggest that the chi-square approximation is very good if $n \geq 3(m^2 + p^2)$ for multivariate normal errors.

Theorems 4 and 5 are useful for relating multivariate tests with the partial F test for multiple linear regression that tests whether a reduced model that omits some of the predictors can be used instead of the full model that uses all p predictors. The partial F test statistic is

$$F_R = \left[\frac{SSE(R) - SSE(F)}{df_R - df_F} \right] / MSE(F)$$

where the residual sums of squares $SSE(F)$ and $SSE(R)$ and degrees of freedom df_F and df_r are for the full and reduced model while the mean square error $MSE(F)$ is for the full model. Let the null hypothesis for the partial F test be $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{0}$ where \mathbf{L} sets the coefficients of the predictors in the full model but not in the reduced model to 0. Seber and Lee (2003, p. 100) shows that

$$F_R = \frac{[\mathbf{L}\hat{\boldsymbol{\beta}}]^T (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1} [\mathbf{L}\hat{\boldsymbol{\beta}}]}{r\hat{\sigma}^2}$$

is distributed as $F_{r,n-p}$ if H_0 is true and the errors are iid $N(0, \sigma^2)$. Note that for multiple linear regression with $m = 1$, $F_R = (n-p)U(\mathbf{L})/r$ since $\hat{\boldsymbol{\Sigma}}_\epsilon^{-1} = 1/\hat{\sigma}^2$. Hence the scaled Hotelling Lawley test statistic is the partial F test statistic extended to $m > 1$ predictor variables by Theorem 4.

By Theorem 5, for example, $rF_R \xrightarrow{D} \chi_r^2$ for a large class of nonnormal error distribution. If $Z_n \sim F_{k,d_n}$, then $Z_n \xrightarrow{D} \chi_k^2/k$ as $d_n \rightarrow \infty$. Hence using the $F_{r,n-p}$ approximation gives a large sample test with correct asymptotic level, and the partial F test is robust

to nonnormality.

Similarly, using an $F_{rm, n-pm}$ approximation for the following test statistics gives large sample tests with correct asymptotic level by Kakizawa (2009) and similar power for large n . The large sample test will have correct asymptotic level as long as the denominator degrees of freedom $d_n \rightarrow \infty$ as $n \rightarrow \infty$, and $d_n = n - pm$ reduces to the partial F test if $m = 1$ and $U(\mathbf{L})$ is used. Then the three test statistics are

$$\frac{-[n - p - 0.5(m - r + 3)]}{rm} \log(\Lambda(\mathbf{L})), \quad \frac{n - p}{rm} V(\mathbf{L}), \quad \text{and} \quad \frac{n - p}{rm} U(\mathbf{L}).$$

Following Khattree and Naik (1999, p. 67) for the Roy's largest root test, if $h = \max(r, m)$, use

$$\frac{n - p - h + r}{h} \lambda_{max}(\mathbf{L}) \approx F(h, n - p - h + r).$$

The simulations in section 5 suggest that this approximation is good for $r = 1$ but poor for $r > 1$. Anderson (1984, p. 333) states that Roy's largest root test has the greatest power if $r = 1$ but is an inferior test for $r > 1$.

Multivariate analogs of tests for multiple linear regression can be derived with appropriate choice of \mathbf{L} . Assume a constant $x_1 = 1$ is in the model. The analog of the ANOVA F test for multiple linear regression is the MANOVA F test that uses $\mathbf{L} = [\mathbf{0} \quad \mathbf{I}_{p-1}]$ to test whether the nontrivial predictors are needed in the model. This test should reject H_0 if the response and residual plots look good, n is large enough and at least one response plot does not look like the corresponding residual plot. A response plot for Y_j will look like a residual plot if the identity line appears almost horizontal, hence the range of \hat{Y}_j is small.

The F_j test of hypotheses uses $\mathbf{L}_j = [0, \dots, 0, 1, 0, \dots, 0]$, where the 1 is in the j th

position, to test whether the j th predictor is needed in the model given that the other $p - 1$ predictors are in the model. This test is an analog of the t tests for multiple linear regression.

The MANOVA partial F test is used to test whether a reduced model is good where the reduced model deletes r of the variables from the full model. For this test, the i th row of \mathbf{L} has a 1 in the position corresponding to the i th variable to be deleted. Omitting the j th variable corresponds to the F_j test while omitting variables x_2, \dots, x_p corresponds to the MANOVA F test. Using $\mathbf{L} = [\mathbf{0} \ \mathbf{I}_k]$ tests whether the last k predictors are needed in the multivariate linear regression model given that the remaining predictors are in the model.

5 Examples and Simulations

The semiparametric and parametric MVN prediction regions applied to the $\hat{\mathbf{z}}_i$ are only conjectured to be large sample prediction regions, but are added to the DD plot as visual aids. Cases below the horizontal line that crosses the identity line correspond to the semiparametric region while cases below the horizontal line that ends at the identity line correspond to the parametric MVN region. A vertical line dropped down from this point of intersection does correspond to a large sample prediction region for multivariate normal data. Note that $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$, and adding a constant $\hat{\mathbf{y}}_f$ to all of the residual vectors does not change the Mahalanobis distances, so the DD plot of the residuals can be used to display the prediction regions.

Example 1. Cook and Weisberg (1999a, p. 351, 433, 447) gives a data set on 82 mussels sampled off the coast of New Zealand. Let $Y_1 = \log(S)$ and $Y_2 = \log(M)$ where S is the shell mass and M is the muscle mass. The predictors are $X_2 = L$, $X_3 = \log(W)$ and $X_4 = H$: the shell length, $\log(\text{width})$ and height. Figures 1 and 2 give the response and residual plots for Y_1 and Y_2 . The response plots show strong linear relationships. For Y_1 , case 79 sticks out while for Y_2 , cases 8, 25 and 48 are not fit well. Highlighted cases had Cook's distance $> \min(0.5, 2p/n)$. See Cook (1977). Figure 3 shows the DD plot of the residual vectors. The plotted points are highly correlated but do not cover the identity line, suggesting an elliptically contoured error distribution that is not multivariate normal. The nonparametric 90% prediction region for the residuals consists of the points to the left of the vertical line $MD = 2.60$. Cases 8, 48 and 79 have especially large distances. The four Hotelling Lawley F_j statistics were greater than 5.77 with pvalues less than 0.005, and the MANOVA F statistic was 337.8 with pvalue ≈ 0 .

The response, residual and DD plots are effective for finding influential cases, for checking linearity and whether the error distribution is multivariate normal or some other elliptically contoured distribution, and for displaying the nonparametric prediction region. Note that cases to the right of the vertical line correspond to cases that are not in their prediction region. These are the cases corresponding to residual vectors with large Mahalanobis distances. Adding a constant does not change the distance, so the DD plot for the residuals is the same as the DD plot for the \hat{z}_i .

Suppose the same model is used except $Y_2 = M$. Then the response and residual plots for Y_1 remain the same, but the plots shown in Figure 4 show curvature about the identity and $r = 0$ lines. Hence the linearity condition is violated. Figure 5 shows that the

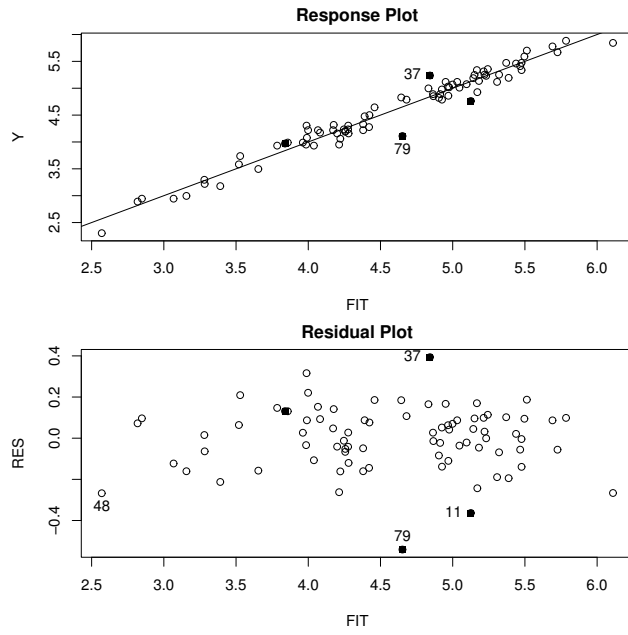


Figure 1: Plots for $Y_1 = \log(S)$.

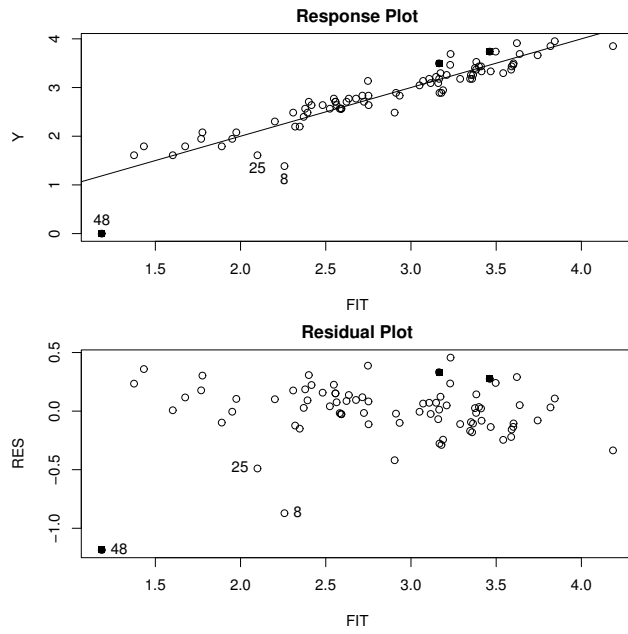


Figure 2: Plots for $Y_2 = \log(M)$.

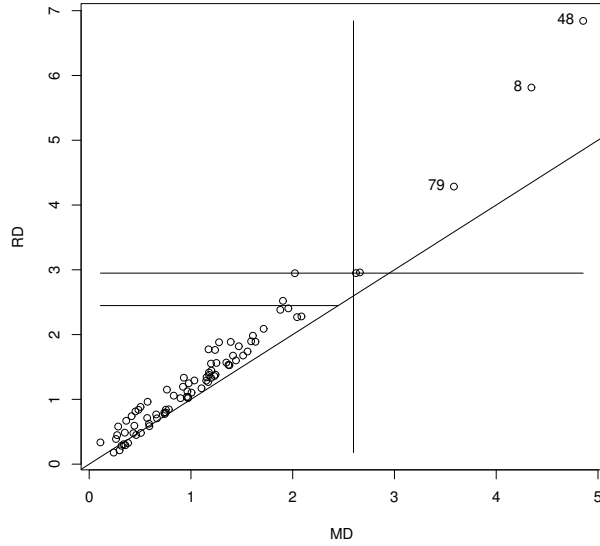


Figure 3: DD Plot of the Residual Vectors for the Mussel Data.

plotted points in the DD plot have correlation well less than one, suggesting that the error distribution is no longer elliptically contoured. The nonparametric 90% prediction region for the residual vectors consists of the points to the left of the vertical line $MD = 2.52$, and contains 95% of the data. Note that the plots can be used to quickly assess whether power transformations have resulted in a linear model, and whether influential cases are present.

Example 2. Consider the one way MANOVA model on the famous iris data set with $n = 150$ and $p = 3$ species of iris: setosa, versicolor, and virginica. The $m = 4$ variables are $Y_1 = \text{sepal length}$, $Y_2 = \text{sepal width}$, $Y_3 = \text{petal length}$ and $Y_4 = \text{petal width}$. See Becker, Chambers and Wilks (1988). Figure 6 shows the response and residual plots for Y_4 . Note that the spread of the three dot plots is similar. The dot plot intersects the identity line at the sample mean of the cases in the dot plot. The setosa cases in

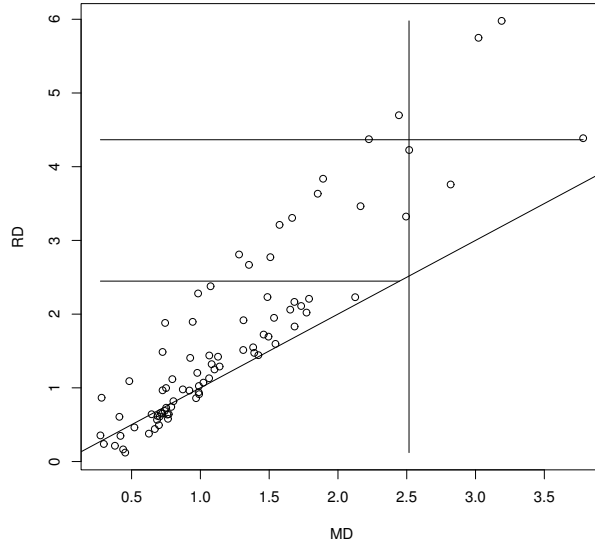


Figure 4: Plot for $Y_2 = M$.

lowest dot plot have a sample mean of 0.246 and the horizontal line $Y_4 = 0.246$ is below the dot plots for versicolor and virginica which have means of 1.326 and 2.026. Hence the mean petal widths differ for the three species, and it is easier to see this difference in the response plot than the residual plot. The plots for the other three variables are similar. Figure 7 shows that the DD plot suggests that the error distribution is elliptically contoured but not multivariate normal. For the one way MANOVA model, a prediction region for \mathbf{y}_f would only be valid for an \mathbf{x}_f which was observed, i.e., for $\mathbf{x}_f = \mathbf{x}_j$, since only observed values of the categorical predictor variables make sense. The 90% nonparametric prediction region corresponds to \mathbf{y} with distances to the left of the vertical line $MD = 3.2$.

A small simulation was used to study the prediction region and the Wilks' Lambda test, the Pillai's trace test, the Hotelling Lawley trace test, and the Roy's largest root

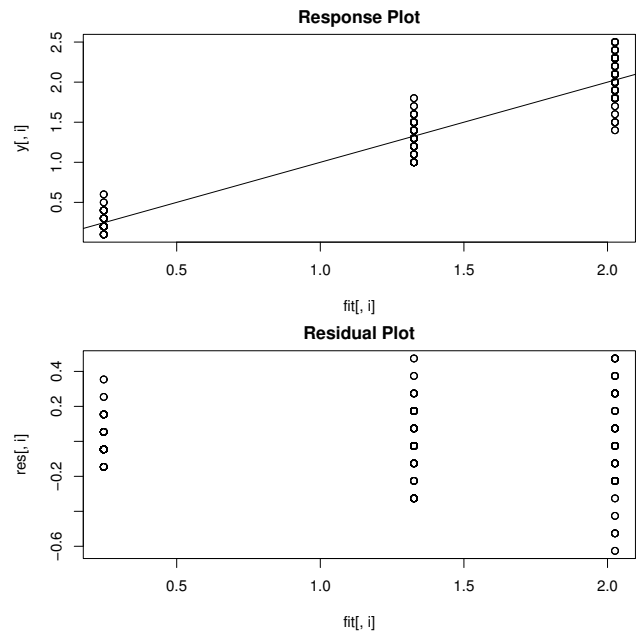


Figure 5: Plots for $Y_4 = \text{Petal Width}$.

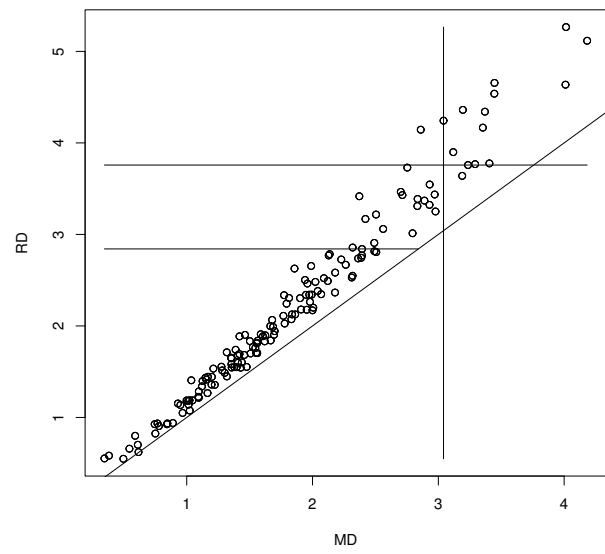


Figure 6: DD Plot of the Residual Vectors for Iris Data.

test for the F_j tests and the MANOVA F test for multivariate linear regression. The first row of \mathbf{B} was always $\mathbf{1}^T$ and the last row of \mathbf{B} was always $\mathbf{0}^T$. When the null hypothesis for the MANOVA F test is true, all but the first row corresponding to the constant are equal to $\mathbf{0}^T$. When $p \geq 3$ and the null hypothesis for the MANOVA F test is false, then the second to last row of \mathbf{B} is $(1, 0, \dots, 0)$, the third to last row is $(1, 1, 0, \dots, 0)$ et cetera as long as the first row is not changed from $\mathbf{1}^T$. First $m \times 1$ error vectors \mathbf{w}_i were generated such that the m errors are iid with variance σ^2 . Let the $m \times m$ matrix $\mathbf{A} = (a_{ij})$ with $a_{ii} = 1$ and $a_{ij} = \psi$ where $0 \leq \psi < 1$ for $i \neq j$. Then $\boldsymbol{\epsilon}_i = \mathbf{A}\mathbf{w}_i$ so that $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}} = \sigma^2 \mathbf{A}\mathbf{A}^T = (\sigma_{ij})$ where the diagonal entries $\sigma_{ii} = \sigma^2[1 + (m - 1)\psi^2]$ and the off diagonal entries $\sigma_{ij} = \sigma^2[2\psi + (m - 2)\psi^2]$ where $\psi = 0.10$. Hence the correlations are $(2\psi + (m - 2)\psi^2)/(1 + (m - 1)\psi^2)$. As ψ gets close to 1, the data clusters about the line in the direction of $(1, \dots, 1)^T$. Used $\mathbf{w}_i \sim N_m(\mathbf{0}, \mathbf{I})$, $\mathbf{w}_i \sim (1 - \tau)N_m(\mathbf{0}, \mathbf{I}) + \tau N_m(\mathbf{0}, 25\mathbf{I})$ with $0 < \tau < 1$ and $\tau = 0.25$ in the simulation, $\mathbf{w}_i \sim$ multivariate t_d with $d = 7$ degrees of freedom, or $\mathbf{w}_i \sim$ lognormal - $E(\text{lognormal})$: where the m components of \mathbf{w}_i were iid with distribution $e^z - E(e^z)$ where $z \sim N(0, 1)$. Only the lognormal distribution is not elliptically contoured.

The simulation used 5000 runs, and H_0 was rejected if the F statistic was greater than $F_{d_1, d_2}(0.95)$ where $P(F_{d_1, d_2} < F_{d_1, d_2}(0.95)) = 0.95$ with $d_1 = rm$ and $d_2 = n - mp$ for the test statistics

$$\frac{-[n - p - 0.5(m - r + 3)]}{rm} \log(\Lambda(\mathbf{L})), \quad \frac{n - p}{rm} V(\mathbf{L}), \quad \text{and} \quad \frac{n - p}{rm} U(\mathbf{L})$$

while $d_1 = h = \max(r, m)$ and $d_2 = n - p - h + r$ for the test statistic

$$\frac{n - p - h + r}{h} \lambda_{\max}(\mathbf{L}).$$

Denote these statistics by W , P , HL and R . Let the coverage be the proportion of times that H_0 is rejected. Want coverage near 0.05 when H_0 is true and coverage close to 1 for good power when H_0 is false. With 5000 runs, coverage outside of (0.04,0.06) suggests that the true coverage is not 0.05. Coverages are tabled for the F_1, F_2, F_{p-1} , and F_p test and for the MANOVA F test denoted by F_M . The null hypothesis H_0 was always true for the F_p test and always false for the F_1 test. When the MANOVA F test was true, H_0 was true for the F_j tests with $j \neq 1$. When the MANOVA F test was false, H_0 was false for the F_j tests with $j \neq p$, but the F_{p-1} test should be hardest to reject for $j \neq p$ by construction of \mathbf{B} and the error vectors.

When the null hypothesis H_0 was true, simulated values started to get close to nominal levels for $n \geq 0.75(m+p)^2$, and were fairly good for $n \geq 1.5(m+p)^2$. The exception was Roy's test which rejects H_0 far too often if $r > 1$. See Table 1 where want values for the F_1 test to be close to 1 since H_0 is false for the F_1 test and want values close to 0.05, otherwise. Roy's test was very good for the F_j tests but very poor for the MANOVA F test. Results are shown for $m = p = 10$. Results from Berndt and Savin (1977) suggest that Pillai's test will reject H_0 less often than Wilks' test which will reject less often than the Hotelling Lawley test.

In Table 2, H_0 is only true for the F_p test where $p = m$, and want values in the F_p column near 0.05. Want values near 1 for high power otherwise. If H_0 is false, often H_0 will be rejected for small n . For example, if $n \geq 10p$, then the response and residual plots should start to look good, and the MANOVA F test should be rejected if there is at least one response variable where the identity line fits the plotted points better than any horizontal line. For the simulated data, had fair power for n not much larger than

Table 1: Test Coverages: MANOVA F H_0 is True.

| \mathbf{w} dist | n | test | F_1 | F_2 | F_{p-1} | F_p | F_M |
|-------------------|-----|------|-------|-------|-----------|-------|-------|
| MVN | 300 | W | 1 | 0.043 | 0.042 | 0.041 | 0.018 |
| MVN | 300 | P | 1 | 0.040 | 0.038 | 0.038 | 0.007 |
| MVN | 300 | HL | 1 | 0.059 | 0.058 | 0.057 | 0.045 |
| MVN | 300 | R | 1 | 0.051 | 0.049 | 0.048 | 0.993 |
| MVN | 600 | W | 1 | 0.048 | 0.043 | 0.043 | 0.034 |
| MVN | 600 | P | 1 | 0.046 | 0.042 | 0.041 | 0.026 |
| MVN | 600 | HL | 1 | 0.055 | 0.052 | 0.050 | 0.052 |
| MVN | 600 | R | 1 | 0.052 | 0.048 | 0.047 | 0.994 |
| MIX | 300 | W | 1 | 0.042 | 0.043 | 0.044 | 0.017 |
| MIX | 300 | P | 1 | 0.039 | 0.040 | 0.042 | 0.008 |
| MIX | 300 | HL | 1 | 0.057 | 0.059 | 0.058 | 0.039 |
| MIX | 300 | R | 1 | 0.050 | 0.050 | 0.051 | 0.993 |
| MVT(7) | 300 | W | 1 | 0.048 | 0.036 | 0.045 | 0.020 |
| MVT(7) | 300 | P | 1 | 0.046 | 0.032 | 0.042 | 0.011 |
| MVT(7) | 300 | HL | 1 | 0.064 | 0.049 | 0.058 | 0.045 |
| MVT(7) | 300 | R | 1 | 0.055 | 0.043 | 0.051 | 0.993 |
| LN | 300 | W | 1 | 0.043 | 0.047 | 0.040 | 0.020 |
| LN | 300 | P | 1 | 0.039 | 0.045 | 0.037 | 0.009 |
| LN | 300 | HL | 1 | 0.057 | 0.061 | 0.058 | 0.041 |
| LN | 300 | R | 1 | 0.049 | 0.055 | 0.050 | 0.994 |

mp. Results are shown for the lognormal distribution.

The same type of data and 5000 runs were used to simulate the prediction regions for \mathbf{y}_f given \mathbf{x}_f for multivariate regression. With $n=100$, $m=2$, and $p=4$, the nominal coverage of the prediction region is 90%, and 92% of the training data is covered. Following Olive (2013), consider the prediction region $\{\mathbf{z} : (\mathbf{z} - T)^T \mathbf{C}^{-1} (\mathbf{z} - T) \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2 \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}} \leq h\}$. Then the ratio of the prediction region volumes

$$\frac{h_i^m \sqrt{\det(\mathbf{C}_i)}}{h_2^m \sqrt{\det(\mathbf{C}_2)}}$$

was recorded where $i = 1$ was the nonparametric region, $i = 2$ was the semiparametric region, and $i = 3$ was the parametric MVN region. Here h_1 and h_2 were the cutoff $D_{(U_n)}(T_i, \mathbf{C}_i)$ for $i = 1, 2$, and $h_3 = \sqrt{\chi_{m,q_n}^2}$.

If, as conjectured, the RMVN estimator is a consistent estimator when applied to the residual vectors instead of iid data, then the volume ratios converge in probability to 1 if the iid zero mean errors $\sim N_m(\mathbf{0}, \Sigma_{\epsilon})$, and the volume ratio converges to 1 for $i = 1$ for a large class of elliptically contoured distributions. These volume ratios were denoted by *voln* and *volm* for the nonparametric and parametric MVN regions. The coverage was the proportion of times the prediction region contained \mathbf{y}_f where *ncov*, *scov* and *mcov* are for the nonparametric, semiparametric and parametric MVN regions.

In the simulations, took $n = 3(m + p)^2$ and $m = p$. Table 3 shows that the coverage of the nonparametric region was close to 0.9 in all cases. The volume ratio *voln* was fairly close to 1 for the three elliptically contoured distributions. Since the volume of the prediction region is proportional to h^m , the volume can be very small if h is too small and m is large. Parametric prediction regions usually give poor estimates of h

Table 2: Test Coverages: MANOVA F H_0 is False.

| n | $m = p$ | test | F_1 | F_2 | F_{p-1} | F_p | F_M |
|-----|---------|------|-------|-------|-----------|-------|-------|
| 30 | 5 | W | 0.012 | 0.222 | 0.058 | 0.000 | 0.006 |
| 30 | 5 | P | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 30 | 5 | HL | 0.382 | 0.694 | 0.322 | 0.007 | 0.579 |
| 30 | 5 | R | 0.799 | 0.871 | 0.549 | 0.047 | 0.997 |
| 50 | 5 | W | 0.984 | 0.955 | 0.644 | 0.017 | 0.963 |
| 50 | 5 | P | 0.971 | 0.940 | 0.598 | 0.012 | 0.871 |
| 50 | 5 | HL | 0.997 | 0.979 | 0.756 | 0.053 | 0.991 |
| 50 | 5 | R | 0.996 | 0.978 | 0.744 | 0.049 | 1 |
| 105 | 10 | W | 0.650 | 0.970 | 0.191 | 0.000 | 0.633 |
| 105 | 10 | P | 0.109 | 0.812 | 0.050 | 0.000 | 0.000 |
| 105 | 10 | HL | 0.964 | 0.997 | 0.428 | 0.000 | 1 |
| 105 | 10 | R | 1 | 1 | 0.892 | 0.052 | 1 |
| 150 | 10 | W | 1 | 1 | 0.948 | 0.032 | 1 |
| 150 | 10 | P | 1 | 1 | 0.941 | 0.025 | 1 |
| 150 | 10 | HL | 1 | 1 | 0.966 | 0.060 | 1 |
| 150 | 10 | R | 1 | 1 | 0.965 | 0.057 | 1 |
| 450 | 20 | W | 1 | 1 | 0.999 | 0.020 | 1 |
| 450 | 20 | P | 1 | 1 | 0.999 | 0.016 | 1 |
| 450 | 20 | HL | 1 | 1 | 0.999 | 0.035 | 1 |
| 450 | 20 | R | 1 | 1 | 0.999 | 0.056 | 1 |

Table 3: Coverages for 90% Prediction Regions.

| \mathbf{w} dist | n | $m = p$ | ncov | scov | mcov | nvol | mvol |
|-------------------|------|---------|-------|-------|-------|-------|-------|
| MVN | 48 | 2 | 0.901 | 0.905 | 0.888 | 0.941 | 0.964 |
| MVN | 300 | 5 | 0.889 | 0.887 | 0.890 | 1.006 | 1.015 |
| MVN | 1200 | 10 | 0.899 | 0.896 | 0.896 | 1.004 | 1.001 |
| MIX | 48 | 2 | 0.912 | 0.927 | 0.710 | 0.872 | 0.097 |
| MIX | 300 | 5 | 0.906 | 0.911 | 0.680 | 0.882 | 0.001 |
| MIX | 1200 | 10 | 0.904 | 0.911 | 0.673 | 0.889 | 0+ |
| MVT(7) | 48 | 2 | 0.903 | 0.910 | 0.825 | 0.914 | 0.646 |
| MVT(7) | 300 | 5 | 0.899 | 0.909 | 0.778 | 0.916 | 0.295 |
| MVT(7) | 1200 | 10 | 0.906 | 0.911 | 0.726 | 0.919 | 0.061 |
| LN | 48 | 2 | 0.912 | 0.926 | 0.651 | 0.729 | 0.090 |
| LN | 300 | 5 | 0.915 | 0.917 | 0.593 | 0.696 | 0.009 |
| LN | 1200 | 10 | 0.912 | 0.916 | 0.593 | 0.679 | 0+ |

when the parametric distribution is misspecified. Hence the parametric MVN region only performed well for multivariate normal data.

6 Conclusions

Multivariate linear regression is a semiparametric method that is nearly as easy to use as multiple linear regression if m is small. The m response and residual plots should be made as well as the DD plot, and the response and residual plots are very useful for the $m = 1$ case of multiple linear regression and experimental design. These plots speed up the model building process for multivariate linear models since the success of power transformations achieving linearity can be quickly assessed and influential cases can be quickly detected. Work is needed on variable selection and on determining the sample sizes for when the tests and prediction regions start to work well. Response and residual plots can look good for $n \geq 10p$, but for testing and prediction regions, may need $n \geq k(m + p)^2$ where $0.5 \leq k \leq 3$ even for well behaved elliptically contoured error distributions.

The *R* software was used to make plots and software. See R Development Core Team (2011). The function `mpredsim` was used to simulate the prediction regions, `mregsim` was used to simulate the tests of hypotheses, and `mregdds` simulated the DD plots for various distributions. The function `mltreg` makes the response and residual plots and computes the F_j , MANOVA F and MANOVA partial F test pvalues while the function `ddplot4` makes the DD plots.

Appendix

Proof of Theorem 1. Let B_n denote the subset of the sample space on which both $\hat{\Sigma}_{1,n}$ and $\hat{\Sigma}_{2,n}$ have inverses. Then $P(B_n) \rightarrow 1$ as $n \rightarrow \infty$. Now

$$\begin{aligned}
D_{\mathbf{x}}^2(\hat{\boldsymbol{\mu}}_j, \hat{\Sigma}_j) &= (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T \hat{\Sigma}_j^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) = \\
&(\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T \left(\frac{\Sigma^{-1}}{a_j} - \frac{\Sigma^{-1}}{a_j} + \hat{\Sigma}_j^{-1} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) = (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T \left(\frac{-\Sigma^{-1}}{a_j} + \hat{\Sigma}_j^{-1} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) + \\
&(\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T \left(\frac{\Sigma^{-1}}{a_j} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) = \frac{1}{a_j} (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T (-\Sigma^{-1} + a_j \hat{\Sigma}_j^{-1}) (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) + \\
&(\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j)^T \left(\frac{\Sigma^{-1}}{a_j} \right) (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j) \\
&= \frac{1}{a_j} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\
&+ \frac{2}{a_j} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j) + \frac{1}{a_j} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j)^T \Sigma^{-1} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j) + \frac{1}{a_j} (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T [a_j \hat{\Sigma}_j^{-1} - \Sigma^{-1}] (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)
\end{aligned}$$

on B_n , and the last three terms are $o_P(1)$ under a) and $O_P(n^{-\delta})$ under b).

Proof of Theorem 2. a) Suppose $(\mathbf{x}_f, \mathbf{y}_f) = (\mathbf{x}_i, \mathbf{y}_i)$. Then

$$D_{\mathbf{y}_i}^2(\hat{\mathbf{y}}_i, \hat{\Sigma}_{\boldsymbol{\epsilon}}) = (\mathbf{y}_i - \hat{\mathbf{y}}_i)^T \hat{\Sigma}_{\boldsymbol{\epsilon}}^{-1} (\mathbf{y}_i - \hat{\mathbf{y}}_i) = \hat{\boldsymbol{\epsilon}}_i^T \hat{\Sigma}_{\boldsymbol{\epsilon}}^{-1} \hat{\boldsymbol{\epsilon}}_i = D_{\hat{\boldsymbol{\epsilon}}_i}^2(\mathbf{0}, \hat{\Sigma}_{\boldsymbol{\epsilon}}).$$

Hence \mathbf{y}_i is in the i th prediction region $\{\mathbf{z} : D_{\mathbf{z}}(\hat{\mathbf{y}}_i, \hat{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{(U_n)}(\hat{\mathbf{y}}_i, \hat{\Sigma}_{\boldsymbol{\epsilon}})\}$ iff $\hat{\boldsymbol{\epsilon}}_i$ is in prediction region $\{\mathbf{z} : D_{\mathbf{z}}(\mathbf{0}, \hat{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{(U_n)}(\mathbf{0}, \hat{\Sigma}_{\boldsymbol{\epsilon}})\}$, but exactly U_n of the $\hat{\boldsymbol{\epsilon}}_i$ are in the latter region by construction, if $D_{(U_n)}$ is unique. Since $D_{(U_n)}$ is the $100(1-\delta)$ th percentile of the D_i asymptotically, $U_n/n \rightarrow 1 - \delta$.

b) Let $P[D_{\mathbf{z}}(E(\mathbf{y}_f), \Sigma_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}(E(\mathbf{y}_f), \Sigma_{\boldsymbol{\epsilon}})] = 1 - \delta$. Since $\Sigma_{\boldsymbol{\epsilon}} > 0$, Theorem 1 shows that if $(\hat{\mathbf{y}}_f, \hat{\Sigma}_{\boldsymbol{\epsilon}}) \xrightarrow{P} (E(\mathbf{y}_f), \Sigma_{\boldsymbol{\epsilon}})$ then $D(\hat{\mathbf{y}}_f, \hat{\Sigma}_{\boldsymbol{\epsilon}}) \xrightarrow{P} D_{\mathbf{z}}(E(\mathbf{y}_f), \Sigma_{\boldsymbol{\epsilon}})$. Hence the percentiles of the distances also converge in probability, and the probability that \mathbf{y}_f is in $\{\mathbf{z} : D_{\mathbf{z}}(\hat{\mathbf{y}}_f, \hat{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}(\hat{\mathbf{y}}_f, \hat{\Sigma}_{\boldsymbol{\epsilon}})\}$ converges to $1 - \delta =$ the probability that \mathbf{y}_f is in $\{\mathbf{z} : D_{\mathbf{z}}(E(\mathbf{y}_f), \Sigma_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}(E(\mathbf{y}_f), \Sigma_{\boldsymbol{\epsilon}})\}$.

c) The asymptotically optimal prediction region is the region with the smallest volume (hence highest density) such that the coverage is $1 - \delta$, as $n \rightarrow \infty$. This region is $\{\mathbf{z} : D_{\mathbf{z}}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}})\}$ if the asymptotically optimal region for the $\boldsymbol{\epsilon}_i$ is $\{\mathbf{z} : D_{\mathbf{z}}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}})\}$. Hence the result follows by b).

Proof of Theorem 3. Multivariate linear regression with least squares satisfies Theorem 2 by Su and Cook (2012). Let (T, \mathbf{C}) be the sample mean and sample covariance matrix (1) applied to the $\hat{\mathbf{z}}_i$. The sample mean and sample covariance matrix of the residual vectors is $(\mathbf{0}, \mathbf{S}_r)$ since least squares was used. Hence the $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$ have sample covariance matrix \mathbf{S}_r , and sample mean $\hat{\mathbf{y}}_f$. Hence $(T, \mathbf{C}) = (\hat{\mathbf{y}}_f, \mathbf{S}_r)$, and the $D_i(\hat{\mathbf{y}}_f, \mathbf{S}_r)$ are used to compute $D_{(U_n)}$.

Proof of Theorem 4. Using the Searle (1982, p. 333) identity $\text{tr}(\mathbf{A}\mathbf{G}^T\mathbf{D}\mathbf{G}\mathbf{C}) = [\text{vec}(\mathbf{G})]^T[\mathbf{C}\mathbf{A} \otimes \mathbf{D}^T][\text{vec}(\mathbf{G})]$, it follows that $(n-p)U(\mathbf{L}) = \text{tr}[\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} \hat{\mathbf{B}}^T \mathbf{L}^T [\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} \mathbf{L} \hat{\mathbf{B}}]$
 $= [\text{vec}(\mathbf{L} \hat{\mathbf{B}})]^T [\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L} \hat{\mathbf{B}})] = T$ where $\mathbf{A} = \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1}$, $\mathbf{G} = \mathbf{L} \hat{\mathbf{B}}$, $\mathbf{D} = [\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1}$, and $\mathbf{C} = \mathbf{I}$. Hence (7) holds.

Proof of Theorem 5. By Su and Cook (2012), $\sqrt{n} \text{vec}(\hat{\mathbf{B}} - \mathbf{B}) \xrightarrow{D} N_{pm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{W})$. Then under H_0 , $\sqrt{n} \text{vec}(\mathbf{L} \hat{\mathbf{B}}) \xrightarrow{D} N_{rm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{L}\mathbf{W}\mathbf{L}^T)$, and $n [\text{vec}(\mathbf{L} \hat{\mathbf{B}})]^T [\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}^{-1} \otimes (\mathbf{L}\mathbf{W}\mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L} \hat{\mathbf{B}})] \xrightarrow{D} \chi_{rm}^2$. This result also holds if \mathbf{W} and $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$ are replaced by $\hat{\mathbf{W}} = n(\mathbf{X}^T \mathbf{X})^{-1}$ and $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}$. Hence under H_0 and using the proof of Theorem 4, $T = (n-p)U(\mathbf{L}) = [\text{vec}(\mathbf{L} \hat{\mathbf{B}})]^T [\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L} \hat{\mathbf{B}})] \xrightarrow{D} \chi_{rm}^2$.

References

Anderson, T.W. (1984), *An Introduction to Multivariate Statistical Analysis*, 2nd ed.,

- Wiley, New York, NY.
- Becker, R.A., Chambers, J.M., and Wilks, A.R. (1988), *The New S Language a Programming Environment for Data Analysis and Graphics*, Wadsworth and Brooks/Cole, Pacific Grove, CA.
- Berndt, E.R. and Savin, N.E. (1977), "Conflict Among Criteria for Testing Hypotheses in the Multivariate Linear Regression Model," *Econometrica*, 45, 1263-1277.
- Cook, R.D. (1977), "Deletion of Influential Observations in Linear Regression," *Technometrics*, 19, 15-18.
- Cook, R.D., and Weisberg, S. (1999a), *Applied Regression Including Computing and Graphics*, Wiley, New York, NY.
- Cook, R.D., and Weisberg, S. (1999b), "Graphs in Statistical Analysis: is the Medium the Message?" *The American Statistician*, 53, 29-37.
- Cook, R.D., and Olive, D.J. (2001), "A Note on Visualizing Response Transformations in Regression," *Technometrics*, 43, 443-449.
- Henderson, H.V., and Searle, S.R. (1977), "Vec and Vech Operators for Matrices, with Some Uses in Jacobians and Multivariate Statistics," *The Canadian Journal of Statistics*, 7, 65-81.
- Johnson, R.A., and Wichern, D.W. (1988), *Applied Multivariate Statistical Analysis*, 2nd ed., Prentice Hall, Englewood Cliffs, NJ.
- Kakizawa, Y. (2009), "Third-Order Power Comparisons for a Class of Tests for Multivariate Linear Hypothesis Under General Distributions," *Journal of Multivariate Analysis*, 100, 473-496.
- Khattree, R., and Naik, D.N. (1999), *Applied Multivariate Statistics with SAS Software*,

- 2nd ed., SAS Institute, Cary, NC.
- Kshirsagar, A.M. (1972), *Multivariate Analysis*, Marcel Dekker, New York, NY.
- Olive, D.J. (2002), “Applications of Robust Distances for Regression,” *Technometrics*, 44, 64-71.
- Olive, D.J. (2007), “Prediction Intervals for Regression Models,” *Computational Statistics and Data Analysis*, 51, 3115-3122.
- Olive, D.J. (2013), “Asymptotically Optimal Regression Prediction Intervals and Prediction Regions for Multivariate Data,” *International Journal of Statistics and Probability*, 2, 90-100.
- Olive, D.J., and Hawkins, D.M. (2005), “Variable Selection for 1D Regression Models,” *Technometrics*, 47, 43-50.
- Olive, D.J., and Hawkins, D.M. (2010), “Robust Multivariate Location and Dispersion,” Preprint, see (www.math.siu.edu/olive/preprints.htm).
- R Development Core Team (2011), “R: a Language and Environment for Statistical Computing,” R Foundation for Statistical Computing, Vienna, Austria, (www.R-project.org).
- Rousseeuw, P.J., and Van Driessen, K. (1999), “A Fast Algorithm for the Minimum Covariance Determinant Estimator,” *Technometrics*, 41, 212-223.
- Searle, S.R. (1982), *Matrix Algebra Useful for Statistics*, Wiley, New York, NY.
- Seber, G.A.F., and Lee, A.J. (2003), *Linear Regression Analysis*, 2nd ed., Wiley, New York, NY.
- Su, Z., and Cook, R.D. (2012), “Inner Envelopes: Efficient Estimation in Multivariate Linear Regression,” *Biometrika*, 99, 687-702.

Zhang, J., Olive, D.J., and Ye, P. (2012), "Robust Covariance Matrix Estimation With Canonical Correlation Analysis," *International Journal of Statistics and Probability*, 1, 119-136.